

Bayesian Hierarchical Models and the Maximum Entropy Principle

Brendon J. Brewer

The University of Auckland

MaxEnt

This is my favourite interpretation of maximum entropy, which I got from Ariel Caticha.

MaxEnt

This is my favourite interpretation of maximum entropy, which I got from Ariel Caticha.

Starting from a prior distribution $\pi(\mathbf{x})$, and given a constraint on probability distributions, the updated distribution is given by maximising the relative entropy

$$H(p; \pi) = - \sum_i p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{\pi(\mathbf{x})} \right) \quad (1)$$

subject to the given constraint and the normalisation condition.

MaxEnt

If you select the right constraint — that the probability of the ‘data’ proposition should now be 1, this is equivalent to the prior \rightarrow posterior update of Bayesian inference, but is more general.

MaxEnt

If you select the right constraint — that the probability of the ‘data’ proposition should now be 1, this is equivalent to the prior \rightarrow posterior update of Bayesian inference, but is more general.

However, the generalisation is hard to use in practice because it is unclear where constraints that refer to probabilities (Jaynes’ “testable information” which is actually untestable) would come from.

Expectation constraints

Consider some function $f(\mathbf{x})$ whose expected value is given (somehow). The updated distribution is given by the well known result:

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp [\lambda f(\mathbf{x})] . \quad (2)$$

The value of λ is selected to ensure that the correct expected value of $f()$ is obtained.

Expectation constraints: general version

With more than one expectation specified:

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp \left[\sum_i \lambda_i f_i(\mathbf{x}) \right]. \quad (3)$$

The values of the λ s are selected to ensure that the correct expected values of $f_i()$ are obtained.

Exponential Example

Consider n positive quantities x_1, \dots, x_n , and an expectation constraint on the average T :

$$\langle T \rangle = \left\langle \frac{1}{n} \sum_{i=1}^n x_i \right\rangle = \mu \quad (4)$$

Exponential Example

Consider n positive quantities x_1, \dots, x_n , and an expectation constraint on the average T :

$$\langle T \rangle = \left\langle \frac{1}{n} \sum_{i=1}^n x_i \right\rangle = \mu \quad (4)$$

The MaxEnt result is independent exponential distributions:

$$p(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\mu} \exp \left(-\frac{x_i}{\mu} \right). \quad (5)$$

Hierarchical Reflex

People often think “but what if μ is really unknown, so this is only the prior conditional on μ ?”. This leads to a hierarchical model with

$$p(\mu, \mathbf{x}) = p(\mu) \prod_{i=1}^n \frac{1}{\mu} \exp\left(-\frac{x_i}{\mu}\right), \quad (6)$$

but the MaxEnt interpretation is lost. Can we bring it back?

Hierarchical Reflex

People often think “but what if μ is really unknown, so this is only the prior conditional on μ ?”. This leads to a hierarchical model with

$$p(\mu, \mathbf{x}) = p(\mu) \prod_{i=1}^n \frac{1}{\mu} \exp\left(-\frac{x_i}{\mu}\right), \quad (6)$$

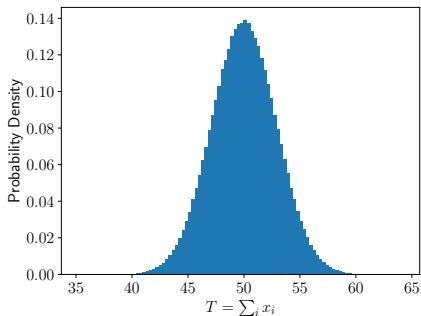
but the MaxEnt interpretation is lost. Can we bring it back?

Answer

Yes, provided we identify the right constraint.

A Different Type of Constraint

Consider a flat prior $\pi(\mathbf{x}) \propto 1$ over a wide domain. This implies a certain prior for the arithmetic mean $T = \frac{1}{n} \sum_i x_i$ that is too narrow and informative (by the central limit theorem). Example with $n = 100$ dimensions and $\text{Uniform}(0, 100)$ prior for each x_i :



Goal

What if we want to control the implied marginal distribution of T , but otherwise keep the distribution as close as possible to π ?

Simple Example — Three-Valued Function

Starting with $\pi(\mathbf{x})$, suppose we want to control the marginal distribution of $f(\mathbf{x})$. For simplicity, assume $f(\mathbf{x})$ can only take values 1, 2, or 3.

Simple Example — Three-Valued Function

Starting with $\pi(\mathbf{x})$, suppose we want to control the marginal distribution of $f(\mathbf{x})$. For simplicity, assume $f(\mathbf{x})$ can only take values 1, 2, or 3. Under any distribution $p(\mathbf{x})$, the probability that $f(\mathbf{x}) = 1$ is

$$P(f(\mathbf{x}) = 1) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 1) d\mathbf{x} \quad (7)$$

Simple Example — Three-Valued Function

Starting with $\pi(\mathbf{x})$, suppose we want to control the marginal distribution of $f(\mathbf{x})$. For simplicity, assume $f(\mathbf{x})$ can only take values 1, 2, or 3. Under any distribution $p(\mathbf{x})$, the probability that $f(\mathbf{x}) = 1$ is

$$P(f(\mathbf{x}) = 1) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 1) d\mathbf{x} \quad (7)$$

Similarly for 2 and 3 we get

$$P(f(\mathbf{x}) = 2) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 2) d\mathbf{x} \quad (8)$$

$$P(f(\mathbf{x}) = 3) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 3) d\mathbf{x} \quad (9)$$

Simple Example — Three-Valued Function

Starting with $\pi(\mathbf{x})$, suppose we want to control the marginal distribution of $f(\mathbf{x})$. For simplicity, assume $f(\mathbf{x})$ can only take values 1, 2, or 3. Under any distribution $p(\mathbf{x})$, the probability that $f(\mathbf{x}) = 1$ is

$$P(f(\mathbf{x}) = 1) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 1) d\mathbf{x} \quad (7)$$

Similarly for 2 and 3 we get

$$P(f(\mathbf{x}) = 2) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 2) d\mathbf{x} \quad (8)$$

$$P(f(\mathbf{x}) = 3) = \int p(\mathbf{x}) \mathbb{1}(f(\mathbf{x}) = 3) d\mathbf{x} \quad (9)$$

Simple Example — Three-Valued Function — Solution

To control the three probabilities (now cast as expected values),
MaxEnt gives the solution

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp [\lambda_1 \mathbb{1}(f(\mathbf{x}) = 1) + \lambda_2 \mathbb{1}(f(\mathbf{x}) = 2) + \lambda_3 \mathbb{1}(f(\mathbf{x}) = 3)]. \quad (10)$$

We tweak λ_1 , λ_2 , and λ_3 to get the desired probabilities.

Simple Example — Three-Valued Function — Solution

To control the three probabilities (now cast as expected values),
MaxEnt gives the solution

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp [\lambda_1 \mathbb{1}(f(\mathbf{x}) = 1) + \lambda_2 \mathbb{1}(f(\mathbf{x}) = 2) + \lambda_3 \mathbb{1}(f(\mathbf{x}) = 3)] . \quad (10)$$

We tweak λ_1 , λ_2 , and λ_3 to get the desired probabilities.

Note that the expression inside the exp is just a funny way of expressing a mapping from f -values to λ values (statisticians may recognise it from regression models involving a factor).

General Solution

From this example, it seems obvious to me (but I haven't proved it) that the general solution to this problem is

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp [g(f(\mathbf{x}))], \quad (11)$$

where the function $g()$ is chosen so that we get the desired marginal distribution for $f(\mathbf{x})$.

General Solution

From this example, it seems obvious to me (but I haven't proved it) that the general solution to this problem is

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) \exp [g(f(\mathbf{x}))], \quad (11)$$

where the function $g()$ is chosen so that we get the desired marginal distribution for $f(\mathbf{x})$.

The exp can be absorbed into the function g if we want, giving a simpler result:

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) g(f(\mathbf{x})). \quad (12)$$

Exponential Example

Returning to the exponential example with a hierarchical prior, the joint prior for the hyperparameter and the parameters is

$$p(\mu, \mathbf{x}) = p(\mu)p(\mathbf{x} | \mu) \quad (13)$$

$$= p(\mu)\mu^{-n} \exp\left(-\frac{\sum_i x_i}{\mu}\right) \quad (14)$$

Exponential Example

Returning to the exponential example with a hierarchical prior, the joint prior for the hyperparameter and the parameters is

$$p(\mu, \mathbf{x}) = p(\mu)p(\mathbf{x} | \mu) \quad (13)$$

$$= p(\mu)\mu^{-n} \exp\left(-\frac{\sum_i x_i}{\mu}\right) \quad (14)$$

and the marginal prior for \mathbf{x} is

$$p(\mathbf{x}) = \int_0^\infty p(\mu)\mu^{-n} \exp\left(-\frac{\sum_i x_i}{\mu}\right) d\mu \quad (15)$$

Exponential Example

The marginal prior for \mathbf{x} is

$$p(\mathbf{x}) = \int_0^\infty p(\mu) \mu^{-n} \exp\left(-\frac{\sum_i x_i}{\mu}\right) d\mu \quad (16)$$

$$p(\mathbf{x}) = \int_0^\infty p(\mu) \mu^{-n} \exp\left(-\frac{nT}{\mu}\right) d\mu \quad (17)$$

which only depends on \mathbf{x} through $T = \frac{1}{n} \sum_i x_i$.

Exponential Example

The marginal prior for \mathbf{x} is

$$p(\mathbf{x}) = \int_0^\infty p(\mu) \mu^{-n} \exp\left(-\frac{\sum_i x_i}{\mu}\right) d\mu \quad (16)$$

$$p(\mathbf{x}) = \int_0^\infty p(\mu) \mu^{-n} \exp\left(-\frac{nT}{\mu}\right) d\mu \quad (17)$$

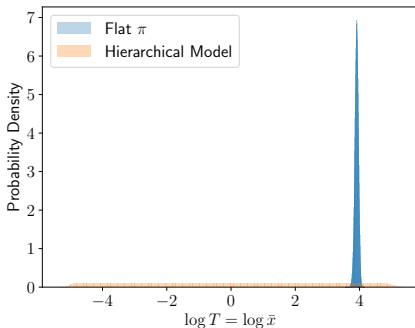
which only depends on \mathbf{x} through $T = \frac{1}{n} \sum_i x_i$. This is therefore a distribution of the form

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) g\left(\frac{1}{n} \sum_i x_i\right) \quad (18)$$

and is MaxEnt with uniform π and a specified marginal distribution for $T = \frac{1}{n} \sum_i x_i$.

What is the Implied Prior on T ?

It is difficult to find the $g(\cdot)$ function for a specified marginal prior, but if we instead just put $\ln(\mu) \sim \text{Uniform}(-5, 5)$ for example, we get something almost log-uniform for $T = \frac{1}{100} \sum_{i=1}^{100} x_i$:



Moral

Moral

Using the hierarchical model is equivalent to using MaxEnt with this specified marginal prior for $T = \frac{1}{n} \sum_i x_i$, fixing the problem with the prior implied by π .

Gaussian Example

Suppose $\pi(\mathbf{x}) \propto 1$ over a wide range. This implies inappropriately narrow priors for $T_1 = \sum_i x_i$ and $T_2 = \sum_i x_i^2$ which may be two quantities of interest.

Gaussian Example

Suppose $\pi(\mathbf{x}) \propto 1$ over a wide range. This implies inappropriately narrow priors for $T_1 = \sum_i x_i$ and $T_2 = \sum_i x_i^2$ which may be two quantities of interest.

Using an extension of the previous result, we can constrain the implied prior for T_1 and T_2 and obtain the MaxEnt distribution:

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) g(f_1(\mathbf{x}), f_2(\mathbf{x})) \quad (19)$$

$$= \pi(\mathbf{x}) g\left(\sum_i x_i, \sum_i x_i^2\right). \quad (20)$$

Gaussian Example

If these were the usual expectation constraints we would have a gaussian distribution over \mathbf{x} , given a μ and a σ . If we let μ and σ be unknown (hierarchical) we get this marginal prior for the \mathbf{x} quantities:

$$p(\mathbf{x}) = \int p(\mu, \sigma) \prod_i \text{Normal}(\mathbf{x}; \mu, \sigma^2) d\mu d\sigma \quad (21)$$

Gaussian Example

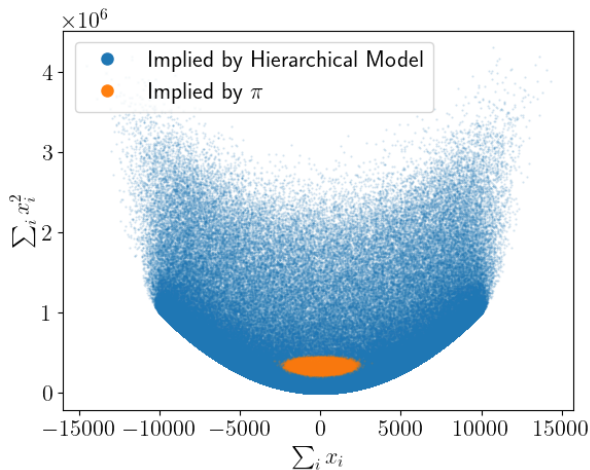
If these were the usual expectation constraints we would have a gaussian distribution over \mathbf{x} , given a μ and a σ . If we let μ and σ be unknown (hierarchical) we get this marginal prior for the \mathbf{x} quantities:

$$p(\mathbf{x}) = \int p(\mu, \sigma) \prod_i \text{Normal}(\mathbf{x}; \mu, \sigma^2) d\mu d\sigma \quad (21)$$

This depends on \mathbf{x} only through the ‘sufficient statistics’ T_1 and T_2 and is thus MaxEnt, of the form

$$p(\mathbf{x}) \propto \pi(\mathbf{x}) g \left(\sum_i x_i, \sum_i x_i^2 \right). \quad (22)$$

The Two Priors



Moral

Moral

Using the hierarchical model is equivalent to using MaxEnt with this specified marginal prior for the sum and sum-of-squares of the x -values, fixing the problem with the prior implied by π .

Conclusions

- At least in some cases, hierarchical models can be thought of as MaxEnt distributions, incorporating a constraint on the marginal distribution of some function of the unknowns.

Conclusions

- At least in some cases, hierarchical models can be thought of as MaxEnt distributions, incorporating a constraint on the marginal distribution of some function of the unknowns.
- The MaxEnt solution for this case is simple with the prior multiplied by a transformed version of the function of interest.

Conclusions

- At least in some cases, hierarchical models can be thought of as MaxEnt distributions, incorporating a constraint on the marginal distribution of some function of the unknowns.
- The MaxEnt solution for this case is simple with the prior multiplied by a transformed version of the function of interest.
- This argument does not rely on any asymptotics, unlike exchangeability/de Finetti justifications for hierarchical models.